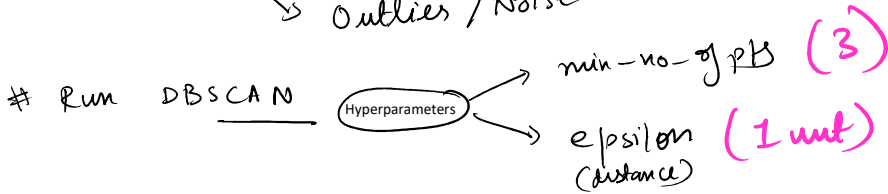
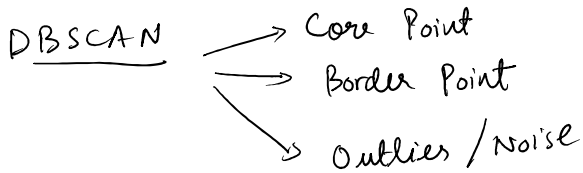
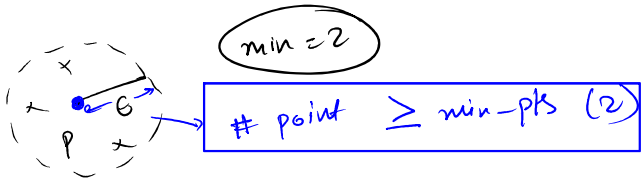
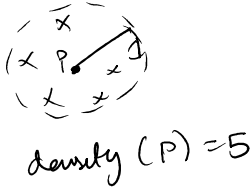
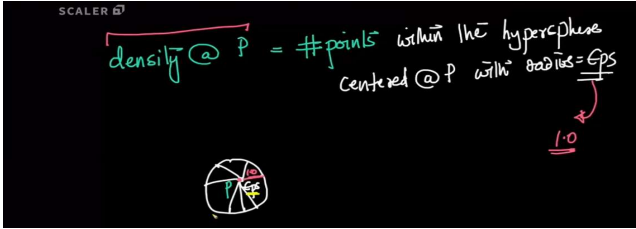


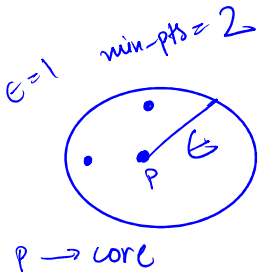
• DBSCAN refers to **Density-based spatial clustering of application with noise**



① Core Point



Border Point \rightarrow ① If P is not a core
② & it lies in the neighbourhood of Q
s.t Q is a core point



Core-point:

{ if P has \geq Minpts in an ϵ ps radius around it then P is a core-pt

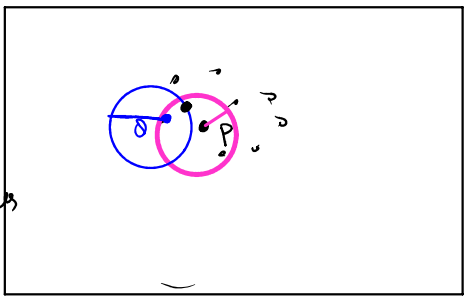
\downarrow ϵ ps, Minpts

P has a dense region around it

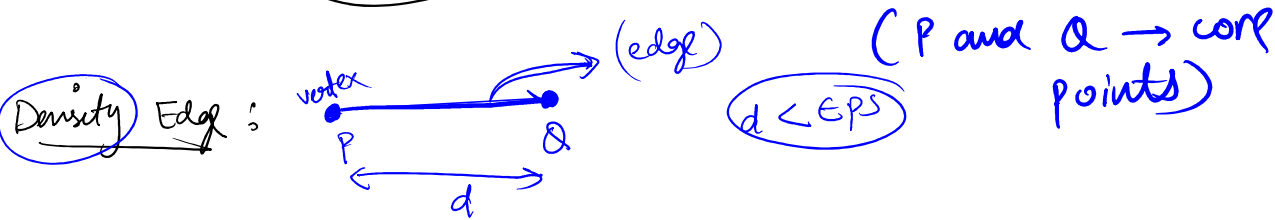
Border Point $\epsilon = 1$
min-pts = 3

P \rightarrow (4) $>$ 3 \Rightarrow P is core
Q \rightarrow (2) $<$ 3 \Rightarrow Q is border

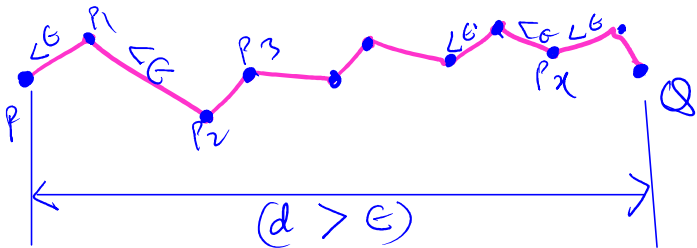
(Q \rightarrow x core but in vicinity of core)



Outlier \rightarrow x core
x Border = outlier



Density connected points = P and Q ($d > \epsilon$ ps but there exists density edge connecting P & Q)



P & Q are core pts

DBSCAN (Algorithm)

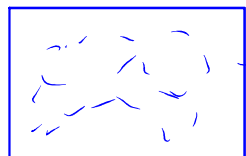
Step 1

($\epsilon = 1$ min pts = 3)

• For each point x_i that belongs to the dataset D , label it as either core point, border point, or noise point.

\rightarrow core

min pts =

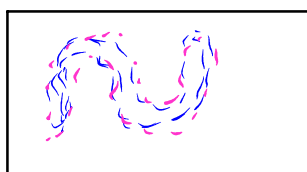


label each pt as

- Core
- Border
- Noise

Step 2

- Remove all the noise points from the dataset



Core

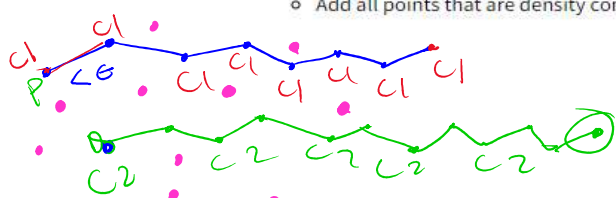
Border

Step 3

Step 3

This step is where things get really interesting. Let's see the 3rd step.

- For each core point P that is not yet assigned to any cluster:
 - create a new cluster with point P
 - Add all points that are density connected to point P , to the P 's cluster



Step 4 → Assign Border points to nearest Cluster

Break till 10:27 PM

Adjusting HypoParams

Industry

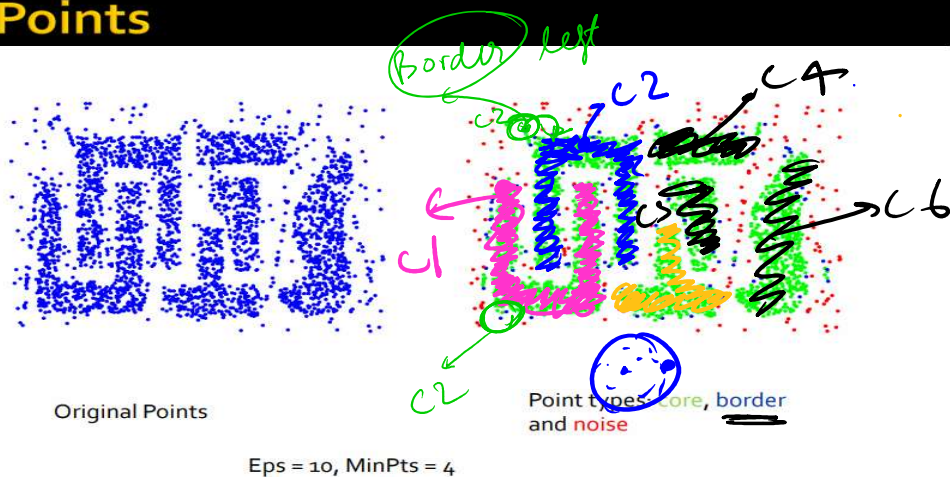
(distance) = Telescopic Search (10³ dimensionality)

min-no-pts $\geq d+1$

libraries $\approx (2d)$

> 20

DBSCAN: Core, Border and Noise Points



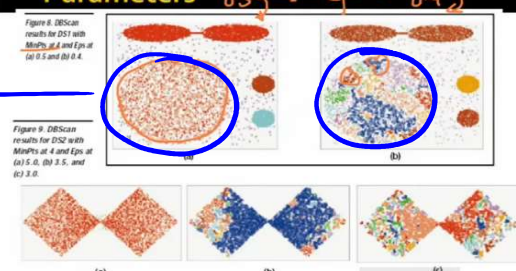
Advantages OF DBSCAN

1. Resistant to noise
2. Handle outliers
3. Non linear data as well
4. Handle different shapes & sizes

DISADVANTAGES

1. Sensitive to parameters
2. Distance Based Algo : Not very good with high dimensions

DBSCAN: Sensitive to Parameters



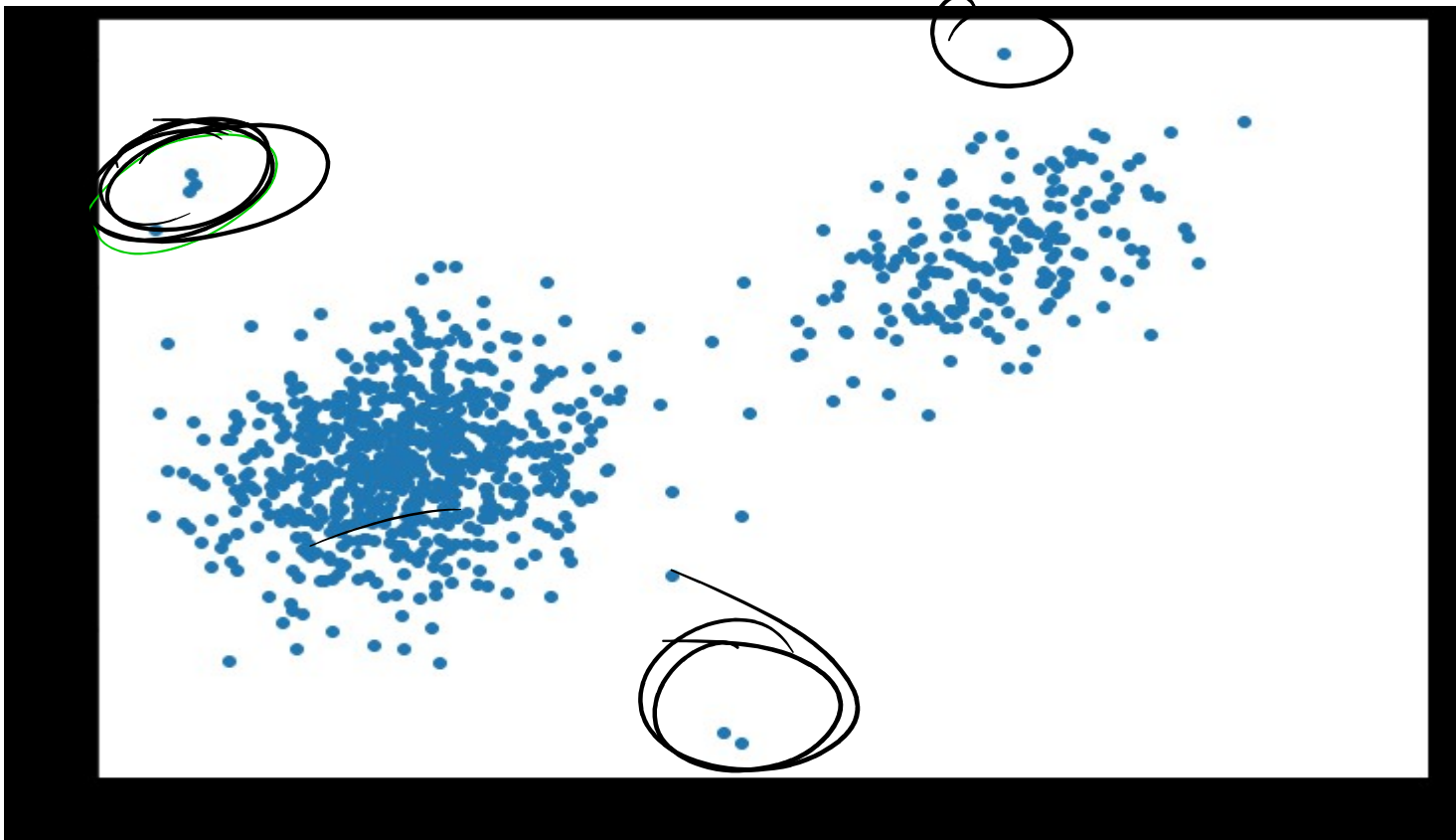
one big cluster

Anomaly

outlier

novelty

Anomaly (outliers)



Lokesh → Categorical Data
↳ Preprocessors

one hot encode
Target encoding

Show term

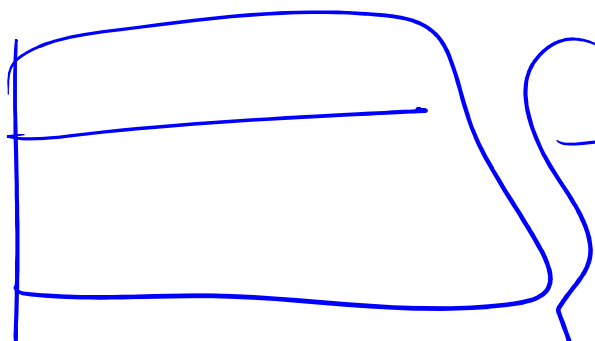
Distance Based Metric

① Meaning of Distance vanishes

Custom Distance Metrics

embeddings

Full DF Gowers
Hamming Distance
Full DF
* (Categorical Data)



one hot encode

model-fit

Target encoding

2)

embedding

model-fit

Manhattan Distance

model-fit

Tools in chest

chest → Harry Dista \rightarrow ~~model~~ \rightarrow df

• $d < \epsilon$
density edge

DBSCAN → Distance + Density Based

$E - M$

posterior probability

Dirichlet Distro → Prior

β

